

[COM4513-6513] Assignment 2: Text Classification with a Feedforward Network

Instructor: Nikos Aletras

The goal of this assignment is to develop a Feedforward network for text classification.

For that purpose, you will implement:

- Text processing methods for transforming raw text data into input vectors for your network (**1 mark**)
- A Feedforward network consisting of:
 - **One-hot** input layer mapping words into an **Embedding weight matrix** (**1 mark**)
 - **One hidden layer** computing the mean embedding vector of all words in input followed by a **ReLU activation function** (**1 mark**)
 - **Output layer** with a **softmax** activation. (**1 mark**)
- The Stochastic Gradient Descent (SGD) algorithm with **back-propagation** to learn the weights of your Neural network. Your algorithm should:
 - Use (and minimise) the **Categorical Cross-entropy loss** function (**1 mark**)
 - Perform a **Forward pass** to compute intermediate outputs (**4 marks**)
 - Perform a **Backward pass** to compute gradients and update all sets of weights (**4 marks**)
 - Implement and use **Dropout** after each hidden layer for regularisation (**2 marks**)
- Discuss how did you choose hyperparameters? You can tune the learning rate (hint: choose small values), embedding size {e.g. 50, 300, 500}, the dropout rate {e.g. 0.2, 0.5} and the learning rate. Please use tables or graphs to show training and validation performance for each hyperparam combination (**2 marks**).
- After training the model, plot the learning process (i.e. training and validation loss in each epoch) using a line plot and report accuracy.
- Re-train your network by using pre-trained embeddings ([GloVe \(https://nlp.stanford.edu/projects/glove/\)](https://nlp.stanford.edu/projects/glove/)) trained on large corpora. Instead of randomly initialising the embedding weights matrix, you should initialise it with the pre-trained weights. During training, you should not update them (i.e. weight freezing) and backprop should stop before computing gradients for updating embedding weights. Report results by performing hyperparameter tuning and plotting the learning process. Do you get better performance? (**3 marks**).
- **BONUS:** Extend you Feedforward network by adding more hidden layers (e.g. one more). How does it affect the performance? Note: You need to repeat hyperparameter tuning, but the number of combinations grows exponentially. Therefore, you need to choose a subset of all possible combinations (+2 extra marks)

Data

The data you will use for Task 2 is a subset of the [AG News Corpus \(http://groups.di.unipi.it/~gulli/AG_corpus_of_news_articles.html\)](http://groups.di.unipi.it/~gulli/AG_corpus_of_news_articles.html) and you can find it in the `./data_topic` folder in CSV format:

- `data_topic/train.csv` : contains 2,400 news articles, 800 for each class to be used for training.
- `data_topic/dev.csv` : contains 150 news articles, 50 for each class to be used for hyperparameter selection and monitoring the training process.
- `data_topic/test.csv` : contains 900 news articles, 300 for each class to be used for testing.

Pre-trained Embeddings

You can download pre-trained GloVe embeddings trained on Common Crawl (840B tokens, 2.2M vocab, cased, 300d vectors, 2.03 GB download) from [here \(http://nlp.stanford.edu/data/glove.840B.300d.zip\)](http://nlp.stanford.edu/data/glove.840B.300d.zip). No need to unzip, the file is large.

Save Memory

To save RAM, when you finish each experiment you can delete the weights of your network using `del w` followed by Python's garbage collector `gc.collect()`

~

```
In [1]: import pandas as pd
import numpy as np
from collections import Counter
import re
import matplotlib.pyplot as plt
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
import random
from time import localtime, strftime
from scipy.stats import spearmanr, pearsonr
import zipfile
import gc

# fixing random seed for reproducibility
random.seed(123)
np.random.seed(123)
```

Transform Raw texts into training and development data

First, you need to load the training, development and test sets from their corresponding CSV files (tip: you can use Pandas dataframes).

In [2]:

In [3]:

Create input representations

To train your Feedforward network, you first need to obtain input representations given a vocabulary. One-hot encoding requires large memory capacity. Therefore, we will instead represent documents as lists of vocabulary indices (each word corresponds to a vocabulary index).

Text Pre-Processing Pipeline

To obtain a vocabulary of words. You should:

- tokenise all texts into a list of unigrams (tip: you can re-use the functions from Assignment 1)
- remove stop words (using the one provided or one of your preference)
- remove unigrams appearing in less than K documents
- use the remaining to create a vocabulary of the top-N most frequent unigrams in the entire corpus.

```
In [4]: stop_words = ['a', 'in', 'on', 'at', 'and', 'or',
                      'to', 'the', 'of', 'an', 'by',
                      'as', 'is', 'was', 'were', 'been', 'be',
                      'are', 'for', 'this', 'that', 'these', 'those', 'you', 'i', 'if',
                      'it', 'he', 'she', 'we', 'they', 'will', 'have', 'has',
                      'do', 'did', 'can', 'could', 'who', 'which', 'what',
                      'but', 'not', 'there', 'no', 'does', 'not', 'so', 've', 'their',
                      'his', 'her', 'they', 'them', 'from', 'with', 'its']
```

Unigram extraction from a document

You first need to implement the `extract_ngrams` function. It takes as input:

- `x_raw` : a string corresponding to the raw text of a document
- `ngram_range` : a tuple of two integers denoting the type of ngrams you want to extract, e.g. (1,2) denotes extracting unigrams and bigrams.
- `token_pattern` : a string to be used within a regular expression to extract all tokens. Note that data is already tokenised so you could opt for a simple white space tokenisation.
- `stop_words` : a list of stop words
- `vocab` : a given vocabulary. It should be used to extract specific features.

and returns:

- a list of all extracted features.

```
In [5]: def extract_ngrams(x_raw, ngram_range=(1,3), token_pattern=r'\b[A-Za-z][A-Za-z]+\b', stop_words=[], vocab=set()):  
  
        return x
```

Create a vocabulary of n-grams

Then the `get_vocab` function will be used to (1) create a vocabulary of ngrams; (2) count the document frequencies of ngrams; (3) their raw frequency. It takes as input:

- `X_raw` : a list of strings each corresponding to the raw text of a document
- `ngram_range` : a tuple of two integers denoting the type of ngrams you want to extract, e.g. (1,2) denotes extracting unigrams and bigrams.
- `token_pattern` : a string to be used within a regular expression to extract all tokens. Note that data is already tokenised so you could opt for a simple white space tokenisation.
- `stop_words` : a list of stop words
- `min_df` : keep ngrams with a minimum document frequency.
- `keep_topN` : keep top-N more frequent ngrams.

and returns:

- `vocab` : a set of the n-grams that will be used as features.
- `df` : a Counter (or dict) that contains ngrams as keys and their corresponding document frequency as values.
- `ngram_counts` : counts of each ngram in vocab

```
In [6]: def get_vocab(X_raw, ngram_range=(1,3), token_pattern=r'\b[A-Za-z][A-Za-z]+\b',  
                     min_df=0, keep_topN=0, stop_words=[]):  
  
        return vocab, df, ngram_counts
```

Now you should use `get_vocab` to create your vocabulary and get document and raw frequencies of unigrams:

```
In [ ]:
```

Then, you need to create vocabulary id -> word and id -> word dictionaries for reference:

In []:

Convert the list of unigrams into a list of vocabulary indices

Storing actual one-hot vectors into memory for all words in the entire data set is prohibitive. Instead, we will store word indices in the vocabulary and look-up the weight matrix. This is equivalent of doing a dot product between an one-hot vector and the weight matrix.

First, represent documents in train, dev and test sets as lists of words in the vocabulary:

In [9]:

In [10]: `x_uni_tr[0]`

```
Out[10]: ['reuters',
          'venezuelans',
          'turned',
          'out',
          'early',
          'large',
          'numbers',
          'sunday',
          'vote',
          'historic',
          'referendum',
          'either',
          'remove',
          'left',
          'wing',
          'president',
          'hugo',
          'chavez',
          'office',
          'give',
          'him',
          'new',
          'mandate',
          'govern',
          'next',
          'two',
          'years']
```

Then convert them into lists of indices in the vocabulary:

In []:

```
In [12]: x_tr[0]
```

```
Out[12]: [5138,  
7303,  
7131,  
2401,  
2311,  
3661,  
3058,  
3338,  
7862,  
6718,  
1632,  
3820,  
6819,  
6518,  
3880,  
256,  
7508,  
1242,  
814,  
6119,  
3920,  
6655,  
2391,  
4338,  
7714,  
3220,  
7982]
```

Put the labels `y` for train, dev and test sets into arrays:

```
In [ ]:
```

Network Architecture

Your network should pass each word index into its corresponding embedding by looking-up on the embedding matrix and then compute the first hidden layer \mathbf{h}_1 :

$$\mathbf{h}_1 = \frac{1}{|x|} \sum_i W_i^e, i \in x$$

where $|x|$ is the number of words in the document and W^e is an embedding matrix $|V| \times d$, $|V|$ is the size of the vocabulary and d the embedding size.

Then \mathbf{h}_1 should be passed through a ReLU activation function:

$$\mathbf{a}_1 = \text{relu}(\mathbf{h}_1)$$

Finally the hidden layer is passed to the output layer:

$$\mathbf{y} = \text{softmax}(\mathbf{a}_1 W^T)$$

where W is a matrix $d \times |\mathcal{Y}|$, $|\mathcal{Y}|$ is the number of classes.

During training, \mathbf{a}_1 should be multiplied with a dropout mask vector (elementwise) for regularisation before it is passed to the output layer.

You can extend to a deeper architecture by passing a hidden layer to another one:

$$\begin{aligned}\mathbf{h}_i &= \mathbf{a}_{i-1} W_i^T \\ \mathbf{a}_i &= \text{relu}(\mathbf{h}_i)\end{aligned}$$

Network Training

First we need to define the parameters of our network by initiliasing the weight matrices. For that purpose, you should implement the `network_weights` function that takes as input:

- `vocab_size` : the size of the vocabulary
- `embedding_dim` : the size of the word embeddings
- `hidden_dim` : a list of the sizes of any subsequent hidden layers (for the Bonus). Empty if there are no hidden layers between the average embedding and the output layer
- `num_clusses` : the number of the classes for the output layer

and returns:

- `W` : a dictionary mapping from layer index (e.g. 0 for the embedding matrix) to the corresponding weight matrix initialised with small random numbers (hint: use `numpy.random.uniform` with from -0.1 to 0.1)

See the examples below for expected outputs. Make sure that the dimensionality of each weight matrix is compatible with the previous and next weight matrix, otherwise you won't be able to perform forward and backward passes. Consider also using `np.float32` precision to save memory.

```
In [15]: def network_weights(vocab_size=1000, embedding_dim=300,
                             hidden_dim=[], num_classes=3, init_val = 0.5):

    return W
```

```
In [15]: W = network_weights(vocab_size=5,embedding_dim=10,hidden_dim=[], num_classes=2)

print('W_emb:', W[0].shape)
print('W_out:', W[1].shape)

W_emb: (5, 10)
W_out: (10, 2)
```

```
In [16]: W = network_weights(vocab_size=3,embedding_dim=4,hidden_dim=[2], num_classes=2)
```

```
In [17]: print('W_emb:', W[0].shape)
print('W_h1:', W[1].shape)
print('W_out:', W[2].shape)

W_emb: (3, 4)
W_h1: (4, 2)
W_out: (2, 2)
```

```
In [18]: W[0]
```

```
Out[18]: array([[ -0.4042875 ,  0.38532683,  0.12724897,  0.22341636],
                [ -0.4838708 ,  0.09443188,  0.05678519, -0.34104034],
                [ -0.3469295 ,  0.19552954, -0.18123357,  0.19197029]],
              dtype=float32)
```

Then you need to develop a `softmax` function (same as in Assignment 1) to be used in the output layer. It takes as input:

- `z` : array of real numbers

and returns:

- `sig` : the softmax of `z`

```
In [19]: def softmax(z):

        return sig
```

Now you need to implement the categorical cross entropy loss by slightly modifying the function from Assignment 1 to depend only on the true label `y` and the class probabilities vector `y_preds` :

```
In [20]: def categorical_loss(y, y_preds):

        return 1
```

```
In [21]: # example for 5 classes

y = 2 #true label
y_preds = softmax(np.array([[ -2.1, 1., 0.9, -1.3, 1.5]]))[0]

print('y_preds: ', y_preds)
print('loss:', categorical_loss(y, y_preds))

y_preds: [0.01217919 0.27035308 0.24462558 0.02710529 0.44573687]
loss: 1.40802648485675
```


Then, implement the `relu` function to introduce non-linearity after each hidden layer of your network (during the forward pass):

$$\text{relu}(z_i) = \max(z_i, 0)$$

and the `relu_derivative` function to compute its derivative (used in the backward pass):

$$\text{relu_derivative}(z_i) = \begin{cases} 0, & \text{if } z_i \leq 0. \\ 1, & \text{otherwise.} \end{cases}$$

Note that both functions take as input a vector z

Hint use `.copy()` to avoid in place changes in array z

```
In [22]: def relu(z):  
  
         return a  
  
         def relu_derivative(z):  
  
         return dz
```

During training you should also apply a dropout mask element-wise after the activation function (i.e. vector of ones with a random percentage set to zero). The `dropout_mask` function takes as input:

- `size` : the size of the vector that we want to apply dropout
- `dropout_rate` : the percentage of elements that will be randomly set to zeros

and returns:

- `dropout_vec` : a vector with binary values (0 or 1)

```
In [23]: def dropout_mask(size, dropout_rate):  
  
         return dropout_vec
```

```
In [24]: print(dropout_mask(10, 0.2))  
print(dropout_mask(10, 0.2))  
  
[1. 1. 0. 1. 1. 1. 1. 1. 0. 1.]  
[1. 1. 1. 1. 0. 1. 1. 0. 1. 1.]
```

Now you need to implement the `forward_pass` function that passes the input `x` through the network up to the output layer for computing the probability for each class using the weight matrices in `w`. The ReLU activation function should be applied on each hidden layer.

- `x` : a list of vocabulary indices each corresponding to a word in the document (input)
- `w` : a list of weight matrices connecting each part of the network, e.g. for a network with a hidden and an output layer: `W[0]` is the weight matrix that connects the input to the first hidden layer, `W[1]` is the weight matrix that connects the hidden layer to the output layer.
- `dropout_rate` : the dropout rate that is used to generate a random dropout mask vector applied after each hidden layer for regularisation.

and returns:

- `out_vals` : a dictionary of output values from each layer: `h` (the vector before the activation function), `a` (the resulting vector after passing `h` from the activation function), its dropout mask vector; and the prediction vector (probability for each class) from the output layer.

```
In [25]: def forward_pass(x, W, dropout_rate=0.2):
```

```
    out_vals = {}

    h_vecs = []
    a_vecs = []
    dropout_vecs = []

    return out_vals
```

```
In [26]: W = network_weights(vocab_size=3,embedding_dim=4,hidden_dim=[5], num_classes=2)
```

```
for i in range(len(W)):
    print('Shape W'+str(i), W[i].shape)

print()
print(forward_pass([2,1], W, dropout_rate=0.5))
```

```
Shape W0 (3, 4)
Shape W1 (4, 5)
Shape W2 (5, 2)
```

```
{'h': [array([-0.04668263, -0.12518334,  0.17532286, -0.32932055], dtype=float32), array([0., 0., 0., 0., 0.])], 'a': [array([0.2286, 0.17532286, 0.], dtype=float32), array([0., 0., 0., 0., 0.])], 'dropout_vec': [array([1., 0., 0., 1.]), array([0., 0., 1., 1., 1.])], 'y': array([0.5, 0.5])}]
```

The `backward_pass` function computes the gradients and update the weights for each matrix in the network from the output to the input. It takes as input

- `x` : a list of vocabulary indices each corresponding to a word in the document (input)
- `y` : the true label
- `W` : a list of weight matrices connecting each part of the network, e.g. for a network with a hidden and an output layer: `W[0]` is the weight matrix that connects the input to the first hidden layer, `W[1]` is the weight matrix that connects the hidden layer to the output layer.
- `out_vals` : a dictionary of output values from a forward pass.
- `learning_rate` : the learning rate for updating the weights.
- `freeze_emb` : boolean value indicating whether the embedding weights will be updated.

and returns:

- `W` : the updated weights of the network.

Hint: the gradients on the output layer are similar to the multiclass logistic regression.

```
In [1]: def backward_pass(x, y, W, out_vals, lr=0.001, freeze_emb=False):  
  
        return W
```

Finally you need to modify SGD to support back-propagation by using the `forward_pass` and `backward_pass` functions.

The `SGD` function takes as input:

- `X_tr` : array of training data (vectors)
- `Y_tr` : labels of `X_tr`
- `W` : the weights of the network (dictionary)
- `X_dev` : array of development (i.e. validation) data (vectors)
- `Y_dev` : labels of `X_dev`
- `lr` : learning rate
- `dropout` : regularisation strength
- `epochs` : number of full passes over the training data
- `tolerance` : stop training if the difference between the current and previous validation loss is smaller than a threshold
- `freeze_emb` : boolean value indicating whether the embedding weights will be updated (to be used by the backward pass function).
- `print_progress` : flag for printing the training progress (train/validation loss)

and returns:

- `weights` : the weights learned
- `training_loss_history` : an array with the average losses of the whole training set after each epoch
- `validation_loss_history` : an array with the average losses of the whole development set after each epoch

```
In [7]: def SGD(X_tr, Y_tr, W, X_dev=[], Y_dev=[], lr=0.001,  
               dropout=0.2, epochs=5, tolerance=0.001, freeze_emb=False, print_progres  
               s=True):  
  
        return W, training_loss_history, validation_loss_history
```

Now you are ready to train and evaluate your neural net. First, you need to define your network using the `network_weights` function followed by SGD with backprop:

```
In [9]: W = network_weights(vocab_size=len(vocab), embedding_dim=300, hidden_dim=[], num_
        classes=3)

        for i in range(len(W)):
            print('Shape W'+str(i), W[i].shape)

        W, loss_tr, dev_loss = SGD(X_tr, Y_tr,
                                   W,
                                   X_dev=X_dev,
                                   Y_dev=Y_dev,
                                   lr=0.001,
                                   dropout=0.2,
                                   freeze_emb=False,
                                   tolerance=0.01,
                                   epochs=100)
```

Plot the learning process:

```
In [ ]:
```

Compute accuracy, precision, recall and F1-Score:

```
In [10]: preds_te = [np.argmax(forward_pass(x, W, dropout_rate=0.0)['y']) for x,y in zip
                    (X_te,Y_te)]
            print('Accuracy:', accuracy_score(Y_te,preds_te))
            print('Precision:', precision_score(Y_te,preds_te,average='macro'))
            print('Recall:', recall_score(Y_te,preds_te,average='macro'))
            print('F1-Score:', f1_score(Y_te,preds_te,average='macro'))
```

Discuss how did you choose model hyperparameters ?

```
In [ ]:
```

Use Pre-trained Embeddings

Now re-train the network using GloVe pre-trained embeddings. You need to modify the `backward_pass` function above to stop computing gradients and updating weights of the embedding matrix.

Use the function below to obtain the embedding matrix for your vocabulary.

```
In [32]: def get_glove_embeddings(f_zip, f_txt, word2id, emb_size=300):

    w_emb = np.zeros((len(word2id), emb_size))

    with zipfile.ZipFile(f_zip) as z:
        with z.open(f_txt) as f:
            for line in f:
                line = line.decode('utf-8')
                word = line.split()[0]

                if word in vocab:
                    emb = np.array(line.strip('\n').split()[1:]).astype(np.float32)
                    w_emb[word2id[word]] += emb

    return w_emb
```

```
In [33]: w_glove = get_glove_embeddings("glove.840B.300d.zip", "glove.840B.300d.txt", word2id)
```

First, initialise the weights of your network using the `network_weights` function. Second, replace the weights of the embedding matrix with `w_glove`. Finally, train the network by freezing the embedding weights:

```
In [ ]:
```

```
In [14]: preds_te = [np.argmax(forward_pass(x, W, dropout_rate=0.0)['y']) for x,y in zip
(X_te,Y_te)]
print('Accuracy:', accuracy_score(Y_te,preds_te))
print('Precision:', precision_score(Y_te,preds_te,average='macro'))
print('Recall:', recall_score(Y_te,preds_te,average='macro'))
print('F1-Score:', f1_score(Y_te,preds_te,average='macro'))
```

Discuss how did you choose model hyperparameters ?

```
In [ ]:
```

Extend to support deeper architectures (Bonus)

Extend the network to support back-propagation for more hidden layers. You need to modify the `backward_pass` function above to compute gradients and update the weights between intermediate hidden layers. Finally, train and evaluate a network with a deeper architecture.

```
In [ ]:
```

```
In [13]: preds_te = [np.argmax(forward_pass(x, W, dropout_rate=0.0)['y']) for x,y in zip
(X_te,Y_te)]
print('Accuracy:', accuracy_score(Y_te,preds_te))
print('Precision:', precision_score(Y_te,preds_te,average='macro'))
print('Recall:', recall_score(Y_te,preds_te,average='macro'))
print('F1-Score:', f1_score(Y_te,preds_te,average='macro'))
```

Full Results

Add your final results here:

Model	Precision	Recall	F1-Score	Accuracy
Average Embedding				
Average Embedding (Pre-trained)				
Average Embedding (Pre-trained) + X hidden layers (BONUS)				

In []: